



清华大学  
Tsinghua University



ICLR 2026

# From Narrow to Panoramic Vision

Attention-Guided Cold-Start Reshapes Multimodal Reasoning

Ruilin Luo, Chufan Shi, Yizhen Zhang, Cheng Yang, Songtao Jiang,  
Tongkun Guan, Ruizhe Chen, Peng Wang, Mingkun Yang, Yujiu Yang, Junyang Lin, Zhibo Yang

Tsinghua University | Qwen Team (Alibaba Group) | University of Southern California

# Table of Contents

01 Background & Motivation

02 Key Discovery 1: VAS

03 Key Discovery 2: Lazy Attention

04 Training-free Interventions

**05 AVAR Framework**

06 Training Details

07 Experimental Results

08 Conclusion

# Background & Motivation

## PROBLEM

Multimodal Large Language Models (MLLMs) often fail to fully leverage visual information in reasoning tasks. Cold-start pretraining shows inconsistent improvements.

## RESEARCH GAP

Understanding how attention works in multimodal reasoning can guide better training strategies. We discovered attention allocation is critical.

## MODEL CATEGORIES



### THE COLD-START QUESTION

Does initializing with image-text pairs actually improve visual reasoning?

# Key Discovery 1: Visual Attention Score (VAS)

## VAS DEFINITION

$$VAS = \frac{\sum \text{Attn}(\text{visual tokens})}{\sum \text{Attn}(\text{system tokens})}$$

Measures attention allocated to visual tokens relative to system tokens

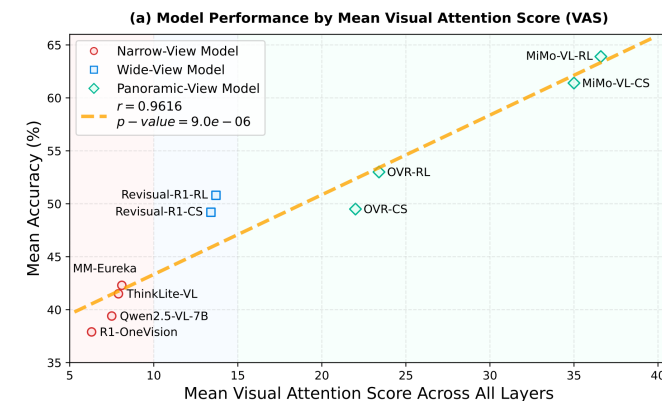
## CORRELATION WITH PERFORMANCE

**Pearson  $r = 0.9616$**

Very strong positive correlation!

## TRAINING INSIGHT

- VAS serves as an effective proxy for measuring multimodal reasoning capability
- Models with higher VAS consistently outperform on math, science, and reasoning benchmarks
- This critical insight guided the entire research direction



# Key Discovery 2: Lazy Attention Localization

Counter-intuitive finding: Multimodal cold-start FAILS to raise VAS

**✗ Multimodal Cold-Start**  
FAILS to raise VAS  
Models don't attend to visual tokens even after visual pretraining

**✓ Text-only Cold-Start**  
Clear VAS increase  
Model learns to attend to text tokens effectively

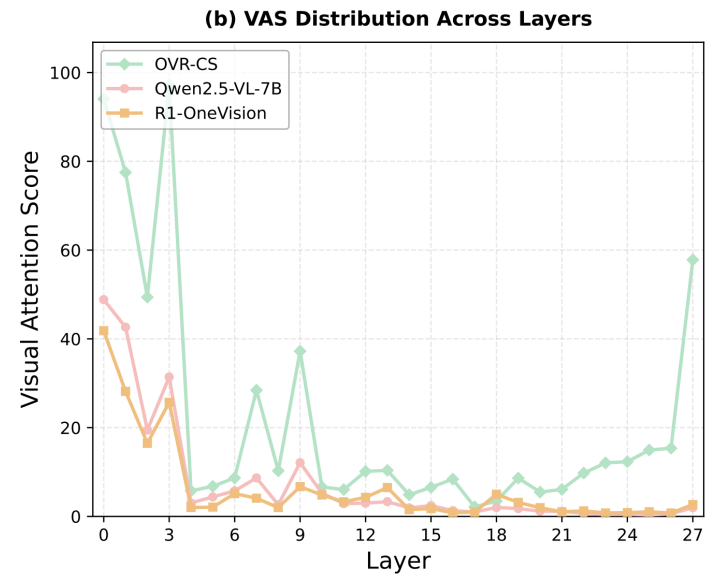
## ROOT CAUSE ANALYSIS

Strong Language Priors

System Token Dominance

"Seen but Not Attended"

**TRAINING IMPLICATION**  
Simply adding visual data in cold-start is not sufficient - this insight led to AVAR framework



# Training-free Interventions for Attention Modulation

Inference-Time Boost for Visual Token Attention in Multimodal Reasoning

## PROBLEM

Even with VAS as a metric and understanding of Lazy Attention, how can we improve attention without model retraining?

## METHOD

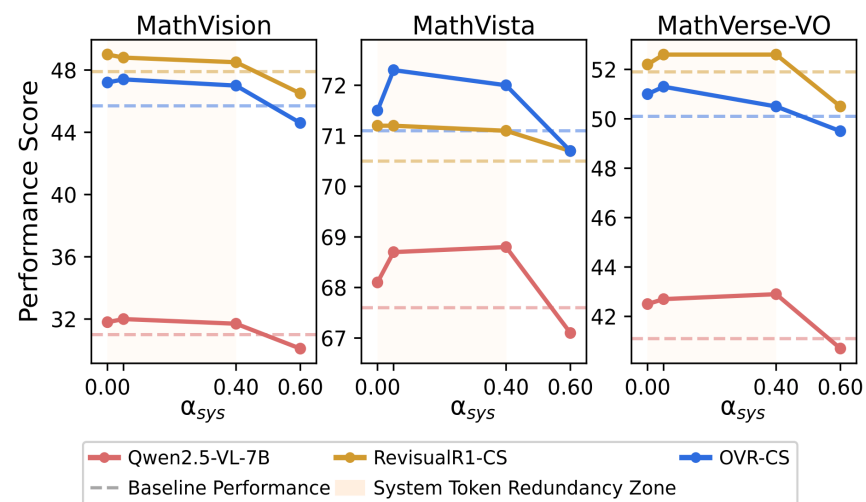
$$Attention\_modified = softmax(Q \cdot K^T + \beta * VisualMask)$$

$\beta$  controls the strength of visual token boost

## SOLUTION

### Attention Modulation at Inference Time

Manipulate attention weights during inference

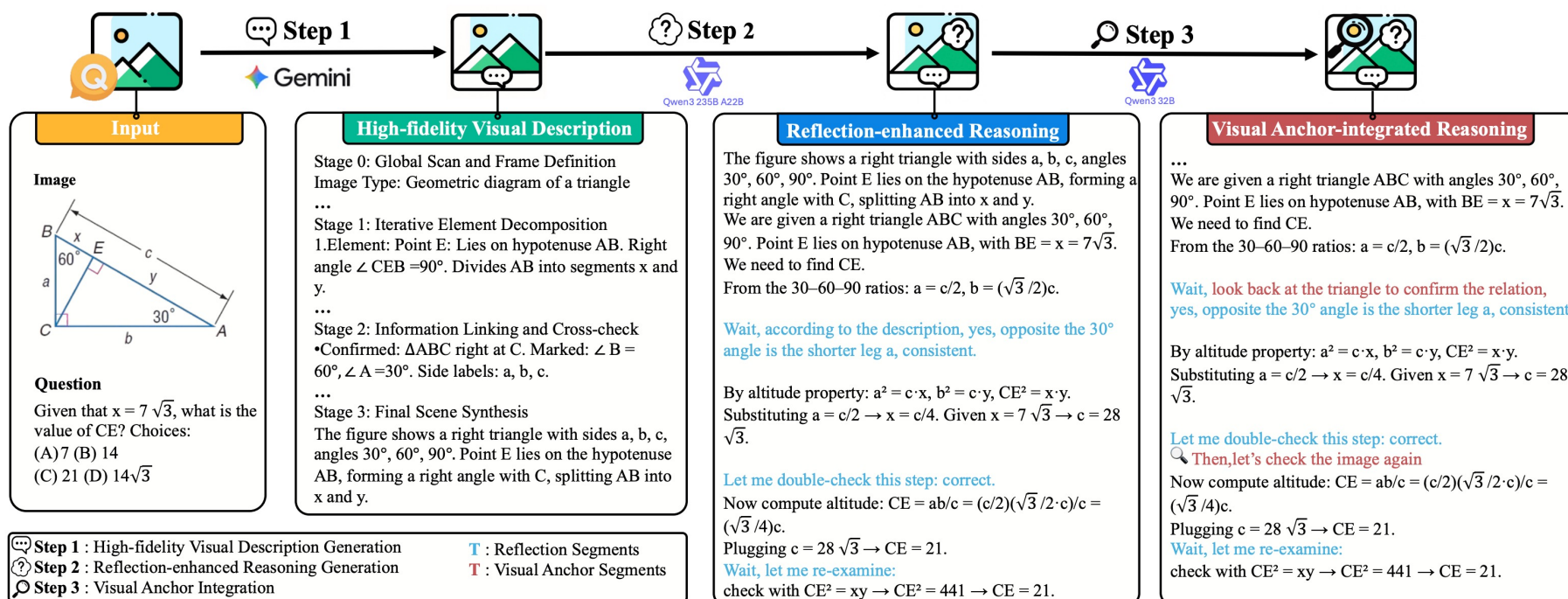


## RESULTS

**+1-2%**

Consistent gains across multiple benchmarks

# AVAR Framework: Attention-Guided Visual Anchoring and Reflection



## KEY RESULT

7.0% average gain across 7 multimodal reasoning benchmarks on Qwen2.5-VL-7B

MathVista | MathVision | MathVerse-VO | MMMU | MMMU-Pro | MMStar | HallusionBench

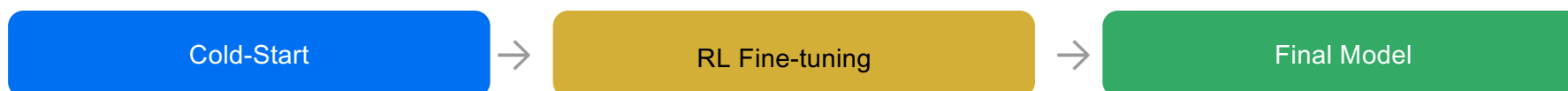
# Results

Model	Math Reasoning			Multidisciplinary		Perception			Avg.
	MathVista	MathVision	MathVerse-VO	MMMU-VAL	MMMU-Pro	MMStar	Hallusion.		
<i>Closed-Source</i>									
GPT-4o	<u>63.8</u>	<u>31.2</u>	-	<u>70.7</u>	<b>54.5</b>	<u>65.1</u>	<u>56.2</u>	-	
Claude-3.7-Sonnet	<b>74.5</b>	<b>58.6</b>	-	<b>75.2</b>	<u>50.1</u>	<b>68.8</b>	<b>58.3</b>	-	
<i>Open-Source General Models</i>									
Qwen2.5-VL-7B	68.2	25.2	41.1	58.1	38.3	62.1	50.7	49.1	
InternVL2.5-8B	64.4	22.0	39.5	56.0	38.2	63.2	51.1	47.8	
LLaVA-OneVision-7B	58.6	18.3	19.3	48.8	35.5	61.7	47.5	41.4	
Llama-3.2-11B-Vision-Instruct	48.6	19.7	18.4	50.7	33.0	49.8	40.3	37.2	
<i>Multimodal Reasoning Models</i>									
Mulberry-7B <sup>†</sup>	63.1	-	42.9	55.0	34.8	61.3	54.1	-	
R1-OneVision	64.1	29.9	40.0	49.1	32.2	52.2	46.0	44.8	
OpenVLThinker	72.3	25.9	44.6	53.0	42.9	59.5	53.0	50.2	
ThinkLite-VL	<b>75.1</b>	<u>32.9</u>	45.8	55.5	40.0	<u>65.0</u>	52.3	<u>53.1</u>	
MM-Eureka-7B	73.0	26.9	48.1	52.0	42.4	<b>65.2</b>	50.7	51.2	
Vision-R1 <sup>†</sup>	73.5	-	47.7	56.3	39.6	64.8	51.9	-	
VLAA-Thinker-7B	68.0	26.4	<u>48.2</u>	55.7	40.9	64.2	50.9	50.6	
Vision-SR1	68.1	26.7	47.1	<u>61.3</u>	<b>43.8</b>	64.1	<u>54.3</u>	52.2	
<i>Our model</i>									
<b>AVAR-Thinker</b>	<b>74.7</b>	<b>37.4</b>	<b>50.4</b>	<b>63.8</b>	<u>42.9</u>	64.1	<b>59.5</b>	<b>56.1</b>	
$\Delta$ over Qwen2.5-VL-7B	+6.5	+12.2	+9.3	+5.7	+4.6	+2.0	+8.8	+7.0	

Configuration	Method Components			Benchmark Performance							
	VARD	AGTO	VARS	MathVista	MathVision	MathVerse-VO	MMStar	MMMU-VAL	MMMU-Pro	Hallusion.	Avg.
Baseline (Qwen2.5-VL-7B)				68.2	25.2	41.1	62.1	58.1	38.3	50.7	49.1
	✓			70.6	32.9	43.5	61.1	55.2	38.7	55.3	51.0
	✓	✓		72.0	34.1	44.0	62.8	58.3	39.8	57.2	52.6
AVAR-Thinker	✓	✓	✓	<b>74.7</b>	<b>37.4</b>	<b>50.4</b>	<b>64.1</b>	<b>63.8</b>	<b>42.9</b>	<b>59.5</b>	<b>56.1</b>

Method	Benchmark Performance							Avg.
	MathVista	MathVision	MathVerse-VO	MMStar	MMMU-val	MMMU-Pro	Hallusion.	
Baseline (Qwen2.5-VL-7B)	68.2	25.2	41.1	<b>62.1</b>	<b>58.1</b>	<u>38.3</u>	50.7	<u>49.1</u>
+ R1-OneVision	63.3	26.3	39.7	54.9	49.9	34.6	43.8	44.6
+ OpenVLThinker	<u>68.9</u>	25.3	37.8	58.7	55.7	36.0	<u>54.1</u>	48.1
+ Vision-SR1	67.6	<u>27.9</u>	<u>42.3</u>	46.9	50.7	36.3	42.1	44.8
+ VARD (Ours)	<b>70.6</b>	<b>32.9</b>	<b>43.5</b>	<u>61.1</u>	<u>55.2</u>	<b>38.7</b>	<b>55.3</b>	<b>51.0</b>

## Two-Phase Training Pipeline



# Conclusion

## VAS

Visual Attention Score  
 $r = 0.9616$  correlation

## Lazy Attention

Root cause of  
cold-start failure

## Training-free

1-2% gains  
without retraining

## AVAR

7.0% gain  
7 benchmarks

## IMPACT

This work provides a fundamental understanding of multimodal attention dynamics and delivers practical solutions that bridge visual pretraining and genuine multimodal reasoning.

# Thank You!

---



清華大學  
Tsinghua University



Qwen



USC University of  
Southern California